

# Probability Theory

Variables: The exact event being tested, but outcome not yet determined  
 $X, Y, Z$  Need to be Random. 3.1

- Sample Outcomes - Any individual potential outcome
- Sample Spaces - The collection of all possible outcomes
- Events - A collection of a subset of  $S$

$\Omega_i$   $\omega_i$   $\omega_1, \omega_2, \dots$   
 $S = \{\omega_1, \omega_2, \omega_3, \dots, \omega_n\}$   
 $A = \{\omega_1, \omega_3, \omega_5, \dots, \omega_n\}$   
 $A \subset S$   $\omega_1, \omega_3, \omega_5, \dots$

	A	1	2	3	4	5	6	7	8	9	10	J	Q	K
♥														
♦														
♣														
♠														

drawing ♥ is an individual outcome  
drawing ♦ is an event  
drawing any card is...

all cards make up the Sample Space

$S = \{A♥, 1♥, 2♥, \dots, K♥, A♦, \dots, K♦, A♣, \dots, A♠, \dots\}$   
 $A = \{A♥, \dots, K♥\}$

$\omega_{A♥} = A♥$

$\omega_i$   
 $P(A♥) = \frac{1}{52} = P(Q♥) = 0.0192 = 1.9\%$

Probability of an outcome:  $P(Q_i) = \frac{1}{52}$  Probability of an event:  $P(A) = \frac{14}{52}$

We can deal with 2 events, A and B  
• Draw 2 cards with Replacement

$P(Q) = \frac{13}{52}$   
 $= \frac{1}{4} = 25\%$   
0.25

	A♥ 1♥	...	Q♥ K♥
A♥	A♥ A♥		A♥ A♥
1♥	A♥ 1♥		Q♥ K♥
...			
Q♥	Q♥ A♥		Q♥ Q♥
K♥	K♥ A♥		K♥ K♥

Size of  $\Omega$  (Norm)

$|\Omega| = 52 \times 52 = 2704$

$P(Q♥ \text{ and } Q♥) = \frac{1}{2704} = 0.00037 = 0.037\%$

$P(♦ \text{ and } ♣) = \frac{13 \times 13}{2704} = \frac{169}{2704} = 0.0625 = 6.25\%$

$P(♦ \text{ and } ♦) = \frac{169}{2704} = 0.0625$

We'll just use decimals

• Draw 2 cards without Replacement

$|\Omega| = 52 \times 51 = 2652$

$P(Q♥ \text{ and } Q♥) = 0$

$P(Q♥ \text{ and } A♥) = \frac{1}{2652} = 0.00038$

$P(♦ \text{ and } ♣) = \frac{13 \times 13}{2652} = \frac{169}{2652} = 0.0637$

$P(♦ \text{ and } ♦) = \frac{13 \times 12}{2652} = \frac{156}{2652} = 0.0588$

Sample Space  $S$   $\Omega$   
Sample Outcome  $\omega_i$   $\omega_i$   
Event  $A$   $A$   
Variable  $X$   $X$

?



# Probability Rules

3.2

$$1. 0 \leq P_r(\omega_i) \leq 1$$

$$2. P_r(\Omega) = 1$$

3. Total Probability of Disjoint (Independent) Events

$$P_r\left(\bigcup_{i=1}^k A_i\right) = \sum_{i=1}^k P_r(A_i)$$

## Independent Events

Conditional Probabilities:  $P(A)$  when you already know  $B$  has happened

$$P(A|B)$$

If  $B$  has no effect on  $A$ , then  $A$  &  $B$  are independent.

$$P(A|B) = P(A)$$

$$P(B|A) = P(B)$$

$$P(A \cap B) = P(A)P(B)$$

- Rolling 2 dice
- Drawing 2 cards w/ replacement
- Flipping a coin

If  $B$  does have an effect on  $A$ , then  $A$  &  $B$  are dependent

- Drawing 2 cards w/out replacement
- Making Purchases: ~~having~~ income

Either way,

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$

Consider a study of people who smoke, relating it to lung cancer

		Cancer		
		$S=1$	$S=0$	
Smoker	$S=1$	0.90	0.01	0.91
	$S=0$	0.03	0.06	0.09
		0.93	0.07	

$$P(S \cap C) = 0.90$$

$$P(S|C) = \frac{0.90}{0.93} = \frac{P(S \cap C)}{P(C)}$$

$$P(C|S) = \frac{0.90}{0.91} = \frac{P(S \cap C)}{P(S)}$$



## Continuous Variables

- Tables work for discrete variables, but not continuous
- We need calculus

Examples include

- Rainfall
- Height  $68'' \rightarrow [67.5'', 68.5'']$
- Time Class ends at 9:20  $\rightarrow [9:18, 9:22]$

We can't do discrete probabilities as exact outcomes are unlikely

$$P(h = 68'' \text{ exactly}) \approx 0$$

So, we use a density function, and we focus on probabilities that an outcome falls within a range.

The probability density function of a variable  $X$  is,

Not a probability.  
Known as a likelihood

$$f_X(\omega) = \frac{1}{\sqrt{2\pi}} e^{-\frac{\omega^2}{2}}$$

p.d.f. can be an number from  $-\infty$  through  $\infty$

$$f_X: \Omega \rightarrow [0, \infty)$$

Graph in Python  
conda install matplotlib  
conda install scipy  
or pip

$$P(X \in A) = \int_A f_X(\omega) d\omega$$

c.d.f. has to be a value between 0 & 1

$$P(X \in A): \Omega \rightarrow [0, 1]$$

$\omega \in A$

$$P(X \in [-1, 0]) = \int_{-1}^0 f_X(\omega) d\omega$$

Area between boundaries of  $A$ .

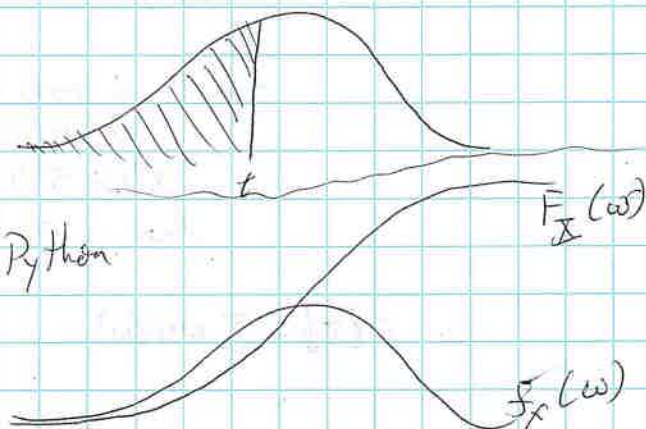
Graph in Desmos



## The Cumulative Density Function

$$F_X(t) = \int_{-\infty}^t f_X(\omega) d\omega$$

Graph in Python





# Random Variables

3.4

A function that maps an outcome to an assigned value

Random Process  $\rightarrow$  outcome  $\rightarrow$  assigned value

A game:

- Roll a 1, Get 5 points
- Roll another odd, Get 1 point
- Roll an even, then lose 2 points

$$\Lambda = \{1, 2, 3, 4, 5, 6\}$$

Sample Space of possible outcomes

$$\Omega = \{5, -2, 1, -2, +1, -2\}$$

Sample Space of possible values

$$X: \Lambda \rightarrow \Omega$$

$$X(1) = 5$$

$$X(2) = -2$$

$$X(3) = 1$$

$$X(4) = -2$$

Do in Python

We often define variables s.t.  $\Omega$  is a set of real values.

In Data Science, we rarely deal with a single value  $\omega \in \mathbb{R}$  but a  $d$ -dimensional vector. Then,

$$\Omega = \{ \omega \mid \omega \in \mathbb{R}^d \}$$

## Expected Values

$E[X]$  is the weighted average of  $\Omega$

• Biased coin:

$$P(H) = 0.4$$

$$P(T) = 0.6$$

$$\Lambda = \{H, T\}$$

• Points:

$$H \rightarrow 5$$

$$T \rightarrow 1$$

$$\Omega = \{5, 1\}$$

$$X(H) = 5$$

$$X(T) = 1$$

• Roll 10 times. On avg., outcomes will be

$$\lambda = H H H H T T T T T T$$

$$X(\omega) = 5 5 5 5 1 1 1 1 1 1$$

$$\text{Average} = \frac{20 + 6}{10} = \frac{26}{10} = 2.6$$

Do in Python

$$E[X] = \sum \omega \cdot P(\omega)$$

$$E[X] = 5 \cdot P(H) + 1 \cdot P(T)$$

$$= 5 \cdot 0.4 + 1 \cdot 0.6$$

$$= 2.0 + 0.6 = 2.6$$



Expected values are linear operations:

$$E[X+Y] = E[X] + E[Y]$$

$$E[cX] = cE[X]$$

Example

Company ABC has 150 employees, each at an avg salary of ~~\$70k~~ \$70k. Bonuses are given based on a 4-tier review:

	R=1	R=2	R=3	R=4
w =	\$1k	\$2k	\$3k	\$4k
P(w) =	0.05	<del>0.05</del> 0.30	0.40	<del>0.30</del> 0.25

As new employee, calculate your expected income

$$E[S] = 70$$

~~PER~~

$$\begin{aligned} E[\text{Pay}] &= E[S + R] = \\ &= E[S] + E[R] \\ &= 70 + (1 \cdot 0.05 + 2 \cdot 0.30 + 3 \cdot 0.40 + 4 \cdot 0.25) \\ &= 70 + (0.05 + 0.60 + 1.20 + 1) \\ &= 70 + 2.85 \\ &= \$72.85k \end{aligned}$$

Variance

Written in terms of the expected value

The variance is a measure of how spread out the data is.

$$\begin{aligned} E[X^2 - 2XE[X] + E[X]^2] \\ E[X^2] - 2E[X]E[X] + E[X]^2 \\ E[X^2] - E[X]^2 \end{aligned}$$

$$\begin{aligned} \text{Var}[X] &= E[(X - E[X])^2] \\ &= E[X^2] - E[X]^2 \end{aligned}$$

$$\sigma_x = \sqrt{\text{Var}[X]} \rightarrow \text{Same units as } X$$

As class Design  
Python programs for  
 $E[X]$   
 $\text{Var}[X]$

Covariance

How much 2 variables vary in relation to each other.

$$\text{Cov}[X, Y] = E[(X - E[X])(Y - E[Y])]$$

In lin alg:

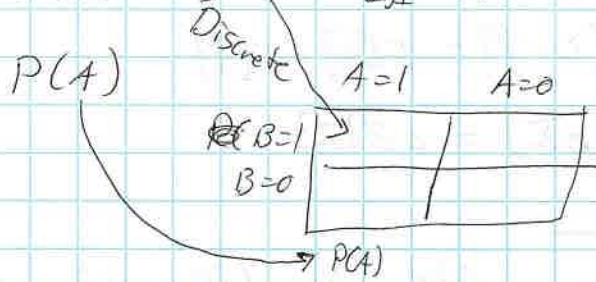
$$\begin{aligned} \bar{X} &= \{x_1 - \bar{x}, x_2 - \bar{x}, \dots\} \\ \text{Cov}[X, Y] &= \frac{\bar{X} \cdot \bar{Y}}{n-1} \end{aligned}$$



Joint Probability:

Continuous  $P(A \cap B) \rightarrow f_{X,Y}(x,y) = P(X=x, Y=y) = P(x \cap y)$

Marginal Probability:



Conditional Probability:  $P(A|B) = \frac{P(A \cap B)}{P(B)} \rightarrow f_{X|Y}(x,y) = \frac{f_{X,Y}(x,y)}{f_Y(y)}$

Discrete Continuous

### Bayes' Rule

You have ~~some~~ models and some data. We want to find the most likely model.  
 (M) (D) (M)

$P(M|D) = \frac{P(D|M) \cdot P(M)}{P(D)}$  ← Sometimes

We can compute this

$P(M|D) = \frac{P(M \cap D)}{P(D)} \leftarrow P(D|M) = \frac{P(M \cap D)}{P(M)} \rightarrow P(M \cap D) = P(D|M) P(M)$

You want to give an award to the best burger in town, based on

- Satisfaction
- Sales

In one day, you find the following sales:

McDonalds: 120 burgers ← ? Great Sales,

Burger's Cafe: 70 burgers

Malt Shoppe: 55 burgers

Polling their customers about satisfaction (come because it's good or because it's convenient):

McDonalds: 40% Satisfaction

Burger's Cafe: 70%

Malt Shoppe: 90%

← ? Great Satisfaction, but low sales.



If you see someone walking around with a burger and satisfied with it, what prob is it that it came from \_\_\_\_\_?

$$P(M_c | S) = \frac{P(S | M_c) \cdot P(M_c)}{P(S)} = \frac{0.40 \cdot \frac{120}{250}}{P(S)} = \frac{0.192}{P(S)}$$

↑ Don't know. That's ok...

$$P(BB | S) = \frac{P(S | B) \cdot P(BB)}{P(S)} = \frac{0.70 \cdot \frac{75}{250}}{P(S)} = \frac{0.21}{P(S)}$$

$$P(MS | S) = \frac{P(S | MS) \cdot P(MS)}{P(S)} = \frac{0.90 \cdot \frac{55}{250}}{P(S)} = \frac{0.198}{P(S)}$$

In Data Science, we use this to find the model <sup>that</sup> most likely fits the data by finding how likely the data fits the model

Take a set of models  $\Omega_M$ . We want  $M \in \Omega_M$  that maximizes  $P(M|D)$

~~$M \in \Omega_M$~~

$$M^* = \underset{M \in \Omega_M}{\operatorname{argmax}} P(M|D)$$

↑  
Known as the  
Posterior  
(The result after  
the test)

~~We want the~~

What is  $\operatorname{argmax}$ ? The argument that maximizes the function.

$$f(x) = \sin(x)$$

$$\max \{f(x)\} = 1$$

$$\underset{x \in [0, \pi]}{\operatorname{argmax}} \{f(x)\} = \pi/2$$

likelihood  
↓

Known as the Prior  
(some knowledge given  
prior to test)  
↓

$$M^* = \underset{M \in \Omega_M}{\operatorname{argmax}} P(M|D) = \underset{M \in \Omega_M}{\operatorname{argmax}} \frac{P(D|M) \cdot P(M)}{P(D)} = \underset{M \in \Omega_M}{\operatorname{argmax}} P(D|M) \cdot P(M)$$

~~$P(D)$~~  is known and fixed, so  $P(D)$  ~~is~~ <sup>is fixed</sup> and doesn't determine the  $\operatorname{argmax}$

Maximum Likelihood Estimator

If  $P(M)$  is a constant, then

$$M^* = \underset{M \in \Omega_M}{\operatorname{argmax}} P(M|D) = \underset{M \in \Omega_M}{\operatorname{argmax}} P(D|M)$$



$P(M)$  a constant? Is it?

In the example, we could have just looked at  $P(S/BB)$ , and it would have given a lot of information.

Multiply by  $P(BB)$  just adds a little more information.

We do want to use  $P(M)$ , but if we don't know it, we can make a first estimate.

What exactly are  $M$  &  $D$ ?

$D = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n\} \in \mathbb{R}^n$  independent observations

$M$ : something that explains structure in the data

Example 1: I want to describe the heights of all students in class

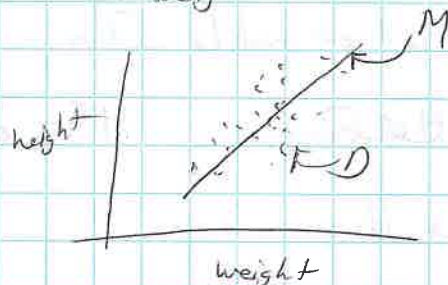
$D = \{x_i\} \in \mathbb{R}$  (one value for each ~~repeated~~ observation)

$M = \bar{x}$

Example 2: I want to compare height with weight

$D = \{\vec{x}_i\} \in \mathbb{R}^2$

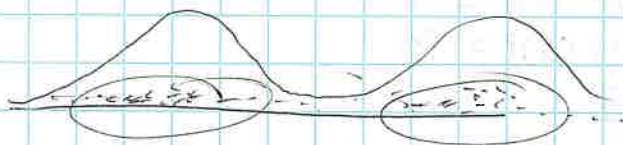
$M$ : Linear Regression



Example 3: Clustering

$D = \{\vec{x}_i\} \in \mathbb{R}^d$

$M$ : Set of points



Example 4: PCA (Principal Component Analysis)

$D = \{\vec{x}_i\} \in \mathbb{R}^d$

$M \in F_k$

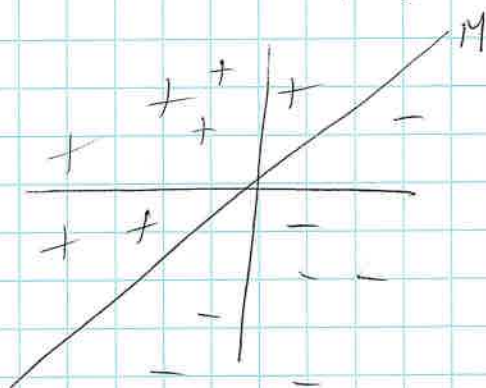
$k$ -dimensional subspace

Dimensionality Reduction



## Example 5: Classification (Linear)

$$D \in \mathbb{R}^{d \times \{-1, +1\}}$$



Let's do an example on Python

- Create an array of random numbers f/normal dist. ( $\mu=0, \sigma=1$ )
- Make histogram
- Find  $\bar{x}$  and  $s$
- Create normal dist from data

$$f_x(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\bar{x})^2}{2\sigma^2}\right)$$

In Python <sup>3080-03.ipynb</sup> <sub>ipynb</sub>

- Plot hist & dist together
- Add normal dist to plot

pdf  $\rightarrow$  like a probability, but not exactly. We can still use it

~~"P"~~

$$P(x_i = x | M = m) \\ \sim P(D_i | M)$$

$$P(D|M) = \prod_i \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

We know the model.  
Given that, what is the probability that  $x$  occurs in this model?



1943

January 2nd

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943

1943